

# Minimizing Condition Number via Convex Programming <sup>\*</sup>

Zhaosong Lu<sup>†</sup>

Ting Kei Pong<sup>‡</sup>

May 11, 2010

Revised: June 22, 2011

## Abstract

In this paper we consider minimizing the spectral condition number of a positive semidefinite matrix over a nonempty closed convex set  $\Omega$ . We show that it can be solved as a convex programming problem, and moreover, the optimal value of the latter problem is achievable. As a consequence, when  $\Omega$  is positive semidefinite representable, it can be cast into a semidefinite programming problem. We then propose a first-order method to solve the convex programming problem. The computational results show that our method is usually faster than the standard interior point solver SeDuMi [16] while producing a comparable solution. We also study a closely related problem, that is, finding an optimal preconditioner for a positive definite matrix. This problem is not convex in general. We propose a convex relaxation for finding positive definite preconditioners. This relaxation turns out to be exact when finding optimal diagonal preconditioners.

**Key words:** condition number, diagonal preconditioner, convex programming, semidefinite programming

**AMS 2000 subject classification:** 90C22, 90C25, 15A12, 65F35

## 1 Introduction

Inspired by Maréchal and Ye [13], we consider the problem of the form

$$\kappa^* = \inf \{ \kappa(X) : X \in \mathcal{S}_+^n \cap \Omega \}, \quad (1)$$

where  $\Omega \subseteq \mathcal{S}^n$  is a nonempty closed convex set,  $\mathcal{S}^n$  is the space of symmetric  $n \times n$  matrices,  $\mathcal{S}_+^n$  is the cone of symmetric positive semidefinite  $n \times n$  matrices, and  $\kappa(X)$  denotes the spectral condition number of  $X$ . We denote by  $\lambda_{\max}(X)$  (resp.  $\lambda_{\min}(X)$ ) the maximal (resp. minimal) eigenvalue of a real symmetric matrix  $X$ . As in [13], for any  $X \in \mathcal{S}_+^n$ , the function  $\kappa$  is defined as

$$\kappa(X) = \begin{cases} \lambda_{\max}(X)/\lambda_{\min}(X) & \text{if } \lambda_{\min}(X) > 0, \\ \infty & \text{if } \lambda_{\min}(X) = 0 \text{ and } \lambda_{\max}(X) > 0, \\ 0 & \text{if } X = 0. \end{cases}$$

---

<sup>\*</sup>This work was supported in part by NSERC Discovery Grant.

<sup>†</sup>Department of Mathematics, Simon Fraser University, Burnaby, BC, V5A 1S6, Canada. (email: [zhaosong@sfu.ca](mailto:zhaosong@sfu.ca)).

<sup>‡</sup>Department of Mathematics, University of Washington, Seattle, Washington 98195, U.S.A. (email: [tkpong@uw.edu](mailto:tkpong@uw.edu)).

It is clear that  $\kappa$  achieves the global minimum value of (1) at 0 if  $0 \in \Omega$ . To avoid some trivial cases, we make the following assumptions on  $\Omega$  in (1) throughout the paper:

**A.1**  $\Omega$  does not contain the zero matrix;

**A.2** The optimal value  $\kappa^*$  of problem (1) is finite.

Problem (1) arises in several applications. For example, Guigues [8] recently applied (1) to estimate the covariance matrix for the Markowitz portfolio selection model (see also [13]). It is easy to show that  $\kappa(\cdot)$  is a quasi-convex function. An approximate solution of (1) can be found by solving a sequence of convex feasibility problems. Indeed, suppose that  $\bar{\kappa}$  and  $\underline{\kappa}$  are the known upper and lower bounds on the optimal value  $\kappa^*$  of (1). Let  $\kappa_l = \underline{\kappa}$ ,  $\kappa_u = \bar{\kappa}$ , and  $v = (\kappa_l + \kappa_u)/2$ . Consider the convex feasibility problem:

$$\text{find } X \in \mathcal{F}_v := \{X \in \mathcal{S}_+^n \cap \Omega, \lambda_{\max}(X) - v\lambda_{\min}(X) \leq 0\}. \quad (2)$$

If  $\mathcal{F}_v = \emptyset$ , we know  $\kappa^* \geq v$  and update  $\kappa_l$  by setting  $\kappa_l \leftarrow v$ . Otherwise,  $\kappa^* \leq v$  and set  $\kappa_u \leftarrow v$ . By repeating this bisection scheme, one can find an  $\epsilon$ -optimal solution of (1) in  $\mathcal{O}(\log \frac{\bar{\kappa} - \underline{\kappa}}{\epsilon})$  number of accesses to the oracle (2) for any given  $\epsilon > 0$ . Though this scheme looks quite simple, it may not be easily implementable as checking whether  $\mathcal{F}_v$  is empty or not can be highly numerically unstable.

Recently, Maréchal and Ye [13] studied problem (1) under the assumption that  $\Omega$  is a compact convex set. They showed that an optimal solution of (1) can be approximated by an exact or an inexact solution of a nonsmooth convex programming problem

$$\min\{\kappa_p(X) : X \in \mathcal{S}_+^n \cap \Omega\}, \quad (3)$$

for some sufficiently large  $p > 0$ , where  $\kappa_p(X) := (\lambda_{\max}(X))^{p+1}/(\lambda_{\min}(X))^p$ . In particular, it is proven in [13] that  $\kappa_p(\cdot)$  is convex for any  $p \geq 0$ , and moreover, every accumulation point of the sequence  $\{X_{p_k}\}$  is an optimal solution of (1) for any  $\{p_k\} \subseteq \mathfrak{R}_+ \rightarrow \infty$ , where  $X_{p_k}$  is an optimal solution of (3) for  $p = p_k$ . It is not known, however, for a given  $\epsilon > 0$ , how large a  $p$  would ensure that an exact or inexact solution of (3) is an  $\epsilon$ -optimal solution of (1).

In this paper we will show that problem (1) can be solved as a convex programming problem, and moreover, the optimal value of the latter problem is achievable. As a consequence, when  $\Omega$  is positive semidefinite representable, it can be cast into a semidefinite programming problem. We then propose a first-order method to solve the convex programming problem and compare its performance with standard interior point (IP) solver SeDuMi for two specific  $\Omega$ 's. The computational results show that our method is usually faster than SeDuMi while producing a comparable solution. We also consider a closely related problem, that is, finding a preconditioner for a positive definite matrix. In particular, assume that  $C \in \mathfrak{R}^{m \times n}$  has full column rank. Consider finding a preconditioner  $X$  for  $C^T C$  so that  $\kappa(X^T C^T C X)$  is minimized, which generally can be formulated as

$$\begin{aligned} \inf \quad & \kappa(X^T C^T C X) \\ \text{s.t.} \quad & X \in \Omega. \end{aligned} \quad (4)$$

Here,  $\Omega$  is a nonempty closed convex subset of  $\mathfrak{R}^{n \times n} \setminus \{0\}$  such that the optimal value of (4) is finite. This problem is typically not convex. We will propose a convex relaxation to it, assuming in addition

that  $X \in \mathcal{S}_+^n$ . Interestingly, this relaxation turns out to be exact when finding diagonal preconditioners. Especially, when  $\Omega$  is a box, finding an optimal diagonal preconditioner can be cast into a cone programming problem.

The rest of the paper is organized as follows. In Section 2, we introduce the notations that are used in this paper and present some basic facts about convex sets for the ease of reference. In Section 3 we consider problem (1) and show that it can be solved as a convex programming problem. In Section 4 we propose a first-order method for solving the convex programming problem and conduct numerical experiments. In Section 5 we study problem (4) and propose a convex relaxation. Finally we present some concluding remarks in Section 6.

## 2 Notations and preliminaries

In this section, we introduce notations used in this paper and present some basic facts about convex sets for the ease of reference.

The symbols  $\mathfrak{R}^n$  and  $\mathfrak{R}_+^n$  denote the  $n$ -dimensional Euclidean space and its nonnegative orthant, respectively. For a vector  $v$ ,  $\|v\|$  denotes the Euclidean norm of  $v$  and  $v_+$  denotes the vector whose  $i$ th entry is  $\max\{v_i, 0\}$ . The  $n$ -dimensional second-order cone will be denoted by  $\mathcal{L}^n$ , that is,

$$\mathcal{L}^n := \left\{ x \in \mathfrak{R}^n : x_1 \geq \sqrt{x_2^2 + \cdots + x_n^2} \right\}.$$

The space of symmetric  $n \times n$  matrices will be denoted by  $\mathcal{S}^n$ . For a matrix  $X \in \mathcal{S}^n$ ,  $\|X\|_F$  denotes the Fröbenius norm of  $X$ ,  $\text{tr}(X)$  denotes the trace of  $X$ , and  $X_{ij}$  denotes the  $(i, j)$ th entry of  $X$ . For matrices  $X, Y \in \mathcal{S}^n$ ,  $\langle X, Y \rangle$  denotes the trace inner product  $\text{tr}(XY)$ , and  $\max\{X, Y\}$  (resp.,  $\min\{X, Y\}$ ) is the matrix whose  $(i, j)$ th entry is  $\max\{X_{ij}, Y_{ij}\}$  (resp.,  $\min\{X_{ij}, Y_{ij}\}$ ). If  $X \in \mathcal{S}^n$  is positive semidefinite (resp., definite), we write  $X \succeq 0$  (resp.,  $X \succ 0$ ). Also, we write  $X \preceq Y$  (resp.,  $X \succeq Y$ ) to mean  $Y - X \succeq 0$  (resp.,  $X - Y \succeq 0$ ). The cone of positive semidefinite (resp., definite) matrices is denoted by  $\mathcal{S}_+^n$  (resp.,  $\mathcal{S}_{++}^n$ ). The cone of nonnegative diagonal  $n \times n$  matrices will be denoted by  $\mathcal{D}_+^n$ . We denote by  $e$  the vector of all ones,  $I$  the identity matrix and  $E$  the matrix of all ones, whose dimensions should be clear from the context. Given a linear operator  $\mathcal{A}$ ,  $\mathfrak{R}(\mathcal{A})$  denotes the range of  $\mathcal{A}$ .

For a set  $S$ , we denote by  $\text{conv}(S)$  and  $\text{cone}(S)$  the convex hull and conical hull of  $S$ , respectively, i.e.,

$$\text{conv}(S) = \left\{ \sum_{i=1}^p \alpha_i s_i : \sum_{i=1}^p \alpha_i = 1, \alpha_i \geq 0, s_i \in S \text{ for all integer } p > 0 \right\}, \quad \text{cone}(S) = \bigcup_{t \geq 0} tS.$$

Notice that if  $S$  is a convex set, then  $\text{cone}(S)$  is a convex set. Given a convex set  $C$ , let  $\text{cl } C$  and  $\text{ri } C$  denote the closure and relative interior of  $C$ , respectively. In the sequel, we will need the following facts about relative interior and closure of convex sets, whose proofs can be found in [14, Chapter 6].

**Proposition 2.1** *Let  $C$  and  $D$  be nonempty convex sets. Then the following statements hold:*

- i) if  $x \in \text{ri } C$  and  $y \in C$ , then  $\alpha x + (1 - \alpha)y \in \text{ri } C$  for any  $\alpha \in (0, 1]$ .*
- ii) if  $\text{ri } C \cap \text{ri } D \neq \emptyset$ , then  $\text{cl } (C \cap D) = \text{cl } C \cap \text{cl } D$  and  $\text{ri } (C \cap D) = \text{ri } C \cap \text{ri } D$ .*

iii) if  $S = \{(1, x) : x \in C\}$ , then  $\text{ri cone}(S) = \{t(1, x) : t > 0, x \in \text{ri } C\}$ .

iv)  $\text{cl}(\text{ri } C) = \text{cl } C$  and  $\text{ri}(\text{cl } C) = \text{ri } C$ .

v)  $\text{ri}(C \times D) = \text{ri } C \times \text{ri } D$ .

In addition, for a closed convex set  $C$ , we denote by  $C_\infty$  the recession cone of  $C$ , that is,

$$C_\infty = \{d : c + td \in C \ \forall t \geq 0\},$$

for any fixed  $c \in C$ . The set is independent of the choice of  $c \in C$  by [1, Proposition 2.1.5] (see also [14, Theorem 8.3]) and is thus well-defined. The following facts about recession cones will be used subsequently, whose proofs can be found in Corollary 2.3.3 and Lemma 2.1.1 of [1], Theorem 8.4 of [14] and Proposition 2.1.11 of [1].

**Proposition 2.2** *Let  $C$  be a nonempty closed convex set. Then the following statements hold.*

i) if  $0 \notin C$ , then  $\text{cl cone}(C) = \text{cone}(C) \cup C_\infty$ .

ii) if  $S = \{(1, x) : x \in C\}$ , then  $\text{cl cone}(S) = \text{cone}(S) \cup \{(0, x) : x \in C_\infty\}$ .

iii)  $C$  is bounded if and only if  $C_\infty = \{0\}$ .

iv) if  $\mathcal{A}$  is a linear map such that  $\mathcal{A}^{-1}(C) \neq \emptyset$ , then  $(\mathcal{A}^{-1}(C))_\infty = \mathcal{A}^{-1}(C_\infty)$ .

Finally, consider the problem of minimizing a real-valued function  $f(x)$  over a certain nonempty feasible region  $\mathcal{F}$  contained in the domain of  $f$  and let  $\bar{f} := \inf\{f(x) : x \in \mathcal{F}\}$ . For  $\epsilon \geq 0$ , we say that  $x_\epsilon$  is an  $\epsilon$ -optimal solution of this problem if  $x_\epsilon \in \mathcal{F}$  and  $f(x_\epsilon) \leq \bar{f} + \epsilon$ .

### 3 Minimizing condition number

In this section we show that problem (1) can be solved as a convex programming problem, and moreover, the optimal value of the latter problem is achievable.

We first show that problem (1) can be reformulated as the following minimization problem with respect to  $(X, t)$ :

$$\lambda^* = \inf \{\lambda_{\max}(X) : X \in t\Omega, t \geq 0, X \succeq I\}. \quad (5)$$

**Theorem 3.1** *The following statements hold:*

i) problem (5) has the same optimal value as (1), that is,  $\lambda^* = \kappa^*$ ;

ii) for any  $\epsilon \geq 0$ , if  $X_\epsilon$  is an  $\epsilon$ -optimal solution of (1), then  $(1/\lambda_{\min}(X_\epsilon), X_\epsilon/\lambda_{\min}(X_\epsilon))$  is an  $\epsilon$ -optimal solution of (5);

iii) for any  $\epsilon \geq 0$ , if  $(t_\epsilon, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (5), then  $X_\epsilon/t_\epsilon$  is an  $\epsilon$ -optimal solution of (1).

*Proof.* By Assumptions A.1 and A.2, we know that  $\mathcal{S}_{++}^n \cap \Omega \neq \emptyset$ . It implies that problem (5) is feasible. Let  $(t, X)$  be a feasible point of (5). Then  $t > 0$  and  $X/t \in \mathcal{S}_{++}^n \cap \Omega$ . Hence,  $X/t$  is a feasible point of (1). Moreover, we have

$$\kappa(X/t) = \kappa(X) = \lambda_{\max}(X)/\lambda_{\min}(X) \leq \lambda_{\max}(X), \quad (6)$$

where the last inequality holds due to  $X \succeq I$ . It then implies  $\kappa^* \leq \lambda^*$ . Now suppose that  $X_\epsilon$  is an  $\epsilon$ -optimal solution of (1) for some  $\epsilon \geq 0$ . Then,  $X_\epsilon \in \mathcal{S}_+^n \cap \Omega$  and  $\kappa(X_\epsilon) \leq \kappa^* + \epsilon$ , which together with Assumptions A.1 and A.2 implies  $\lambda_{\min}(X_\epsilon) > 0$ . It is then straightforward to verify that  $(1/\lambda_{\min}(X_\epsilon), X_\epsilon/\lambda_{\min}(X_\epsilon))$  is a feasible point of (5). Furthermore, we have

$$\lambda^* \leq \lambda_{\max}(X_\epsilon/\lambda_{\min}(X_\epsilon)) = \kappa(X_\epsilon) \leq \kappa^* + \epsilon. \quad (7)$$

Due to the arbitrariness of  $\epsilon$ , we conclude that  $\lambda^* \leq \kappa^*$ . Thus, we have  $\lambda^* = \kappa^*$ , and so statement (i) holds. Moreover, it follows from (7) and statement (i) that  $(1/\lambda_{\min}(X_\epsilon), X_\epsilon/\lambda_{\min}(X_\epsilon))$  is an  $\epsilon$ -optimal solution of (5), and hence statement (ii) holds. Next we show that statement (iii) holds. Indeed, suppose  $(t_\epsilon, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (5) for some  $\epsilon \geq 0$ . We see that  $t_\epsilon > 0$  and  $X_\epsilon/t_\epsilon$  is a feasible point of (1). Replacing  $t$  and  $X$  by  $t_\epsilon$  and  $X_\epsilon$ , respectively in (6), and using statement (i), we obtain that

$$\kappa(X_\epsilon/t_\epsilon) \leq \lambda_{\max}(X_\epsilon) \leq \lambda^* + \epsilon = \kappa^* + \epsilon,$$

which implies that  $X_\epsilon/t_\epsilon$  is an  $\epsilon$ -optimal solution of (1).  $\blacksquare$

We see from Theorem 3.1 that problem (1) can be solved as (5). Notice that the objective function and the feasible region of (5), denoted by  $\mathcal{F}$ , are convex. Nevertheless,  $\mathcal{F}$  generally is not closed. For example, let  $\Omega = \{X \in \mathcal{S}^n : X \succeq I\}$ . Then we see that  $\{(t_k, X_k)\} = \{(1/k, I)\} \subseteq \mathcal{F}$  but  $(t_k, X_k) \rightarrow (0, I) \notin \mathcal{F}$ , which implies that  $\mathcal{F}$  is not closed. We next provide a necessary and sufficient condition for the closedness of  $\mathcal{F}$ .

**Theorem 3.2** *The feasible region  $\mathcal{F}$  of (5) is closed if and only if  $\mathcal{S}_{++}^n \cap \Omega_\infty = \emptyset$ .*

*Proof.* By Assumption A.2 and Theorem 3.1 (i), we know that the optimal value of (5) is finite, which implies that  $\mathcal{S}_{++}^n \cap \Omega \neq \emptyset$ . Let  $X \in \text{ri } \Omega$  and  $Y \in \mathcal{S}_{++}^n \cap \Omega$ . It follows from Proposition 2.1 (i) that  $\{\alpha X + (1 - \alpha)Y : \alpha \in (0, 1]\} \subseteq \text{ri } \Omega$ . Hence we have that  $\alpha X + (1 - \alpha)Y \in \mathcal{S}_{++}^n \cap \text{ri } \Omega$  when  $0 < \alpha \ll 1$ . Thus  $\mathcal{S}_{++}^n \cap \text{ri } \Omega \neq \emptyset$ . We now define

$$\mathcal{K} = \{t(1, X) : t \geq 0, X \in \Omega\}, \quad \tilde{\mathcal{K}} = \{(t, X) : t \in \mathfrak{R}, X \succeq I\}. \quad (8)$$

It follows from Proposition 2.1 (iii) and (v) that

$$\text{ri } \mathcal{K} = \{t(1, X) : t > 0, X \in \text{ri } \Omega\}, \quad \text{ri } \tilde{\mathcal{K}} = \{(t, X) : t \in \mathfrak{R}, X \succ I\}. \quad (9)$$

Since  $\mathcal{S}_{++}^n \cap \text{ri } \Omega \neq \emptyset$ , we see that  $\text{ri } \mathcal{K} \cap \text{ri } \tilde{\mathcal{K}} \neq \emptyset$ . Using this result, and Propositions 2.1 (ii) and 2.2 (ii), we obtain that

$$\begin{aligned} \text{cl } (\mathcal{K} \cap \tilde{\mathcal{K}}) &= \text{cl } \mathcal{K} \cap \text{cl } \tilde{\mathcal{K}} = (\mathcal{K} \cup \{(0, X) : X \in \Omega_\infty\}) \cap \tilde{\mathcal{K}}, \\ &= (\mathcal{K} \cap \tilde{\mathcal{K}}) \cup \{(0, X) : X \in \Omega_\infty, X \succeq I\}, \end{aligned}$$

which together with the definitions of  $\mathcal{K}$  and  $\tilde{\mathcal{K}}$  implies that  $\mathcal{F} = \mathcal{K} \cap \tilde{\mathcal{K}}$  is closed if and only if  $\{(0, X) : X \in \Omega_\infty, X \succeq I\} \subseteq \mathcal{K}$ . Further, by the definition of  $\mathcal{K}$ , the latter condition holds if and only if  $\{(0, X) : X \in \Omega_\infty, X \succeq I\} = \emptyset$ , or equivalently,  $\mathcal{S}_{++}^n \cap \Omega_\infty = \emptyset$ . Thus the conclusion holds.  $\blacksquare$

We observe from Theorems 3.1 and 3.2 that when  $\mathcal{S}_{++}^n \cap \Omega_\infty = \emptyset$ , problem (1) can be solved as convex programming problem (5). In particular, for the case where  $\Omega$  is a compact convex set that is assumed in [13],  $\mathcal{S}_{++}^n \cap \Omega_\infty = \emptyset$  holds since  $\Omega_\infty = \{0\}$  by Proposition 2.2 (iii). Given that  $\mathcal{S}_{++}^n \cap \Omega_\infty = \emptyset$  generally may not hold, we will further consider a relaxation of (5):

$$\mu^* = \inf \{ \lambda_{\max}(X) : (t, X) \in \Xi, X \succeq I \}, \quad (10)$$

where

$$\Xi := \{t(1, X) : t \geq 0, X \in \Omega\} \cup \{(0, X) : X \in \Omega_\infty\}. \quad (11)$$

In view of (8) and Proposition 2.2 (ii), we see that  $\Xi = \text{cl } \mathcal{K}$ , and hence  $\Xi$  is closed and convex. We next show that problem (1) can be solved as convex programming problem (10).

**Theorem 3.3** *The following statements hold:*

- i) *problem (10) has the same optimal value as (1), that is,  $\mu^* = \kappa^*$ ;*
- ii) *for any  $\epsilon \geq 0$ , if  $X_\epsilon$  is an  $\epsilon$ -optimal solution of (1), then  $(1/\lambda_{\min}(X_\epsilon), X_\epsilon/\lambda_{\min}(X_\epsilon))$  is an  $\epsilon$ -optimal solution of (10);*
- iii) *for any  $\epsilon > 0$ , if  $(t_\epsilon, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (10) for some  $t_\epsilon > 0$ , then  $X_\epsilon/t_\epsilon$  is an  $\epsilon$ -optimal solution of (1);*
- iv) *for any  $\epsilon > 0$ , if  $(0, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (10), then  $\bar{X} + \alpha X_\epsilon$  is a  $2\epsilon$ -optimal solution of (1), provided that  $\alpha \geq \underline{\alpha} := \max\{[\lambda_{\max}(\bar{X}) - (\mu^* + 2\epsilon)\lambda_{\min}(\bar{X})]/\epsilon, 1 - \lambda_{\min}(\bar{X}), 0\}$ , where  $\bar{X}$  is an arbitrary point in  $\Omega$ .*

*Proof.* In view of (5) and (10), we see that  $\mu^* \leq \lambda^*$ , which together with Theorem 3.1 (i) implies that  $\mu^* \leq \kappa^*$ . We now simultaneously show that  $\kappa^* \leq \mu^*$  and statements (iii) and (iv) hold. For any  $\epsilon > 0$ , suppose that  $(t_\epsilon, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (10). We now consider two cases:  $t_\epsilon > 0$  or  $t_\epsilon = 0$ . First, assume that  $t_\epsilon > 0$ . We observe that  $\lambda_{\max}(X_\epsilon) \leq \mu^* + \epsilon \leq \lambda^* + \epsilon$ . Thus,  $(t_\epsilon, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (5). It then follows from Theorem 3.1 (iii) that  $X_\epsilon/t_\epsilon$  is an  $\epsilon$ -optimal solution of (1), i.e., statement (iii) holds. Moreover, we have

$$\kappa^* \leq \kappa(X_\epsilon/t_\epsilon) = \kappa(X_\epsilon) \leq \lambda_{\max}(X_\epsilon) \leq \mu^* + \epsilon, \quad (12)$$

where the second inequality follows from  $X_\epsilon \succeq I$ . Next, assume that  $t_\epsilon = 0$ . We can observe that

$$I \preceq X_\epsilon \in \Omega_\infty, \quad \lambda_{\max}(X_\epsilon) \leq \mu^* + \epsilon. \quad (13)$$

Let  $\bar{X}$  be an arbitrary point in  $\Omega$  and  $\underline{\alpha}$  be defined above. In view of (13) and the definition of  $\underline{\alpha}$ , it follows that when  $\alpha \geq \underline{\alpha}$ , we have  $\bar{X} + \alpha X_\epsilon \in \mathcal{S}_{++}^n \cap \Omega$  and

$$\frac{\lambda_{\max}(\bar{X}) + \alpha \lambda_{\max}(X_\epsilon)}{\lambda_{\min}(\bar{X}) + \alpha} \leq \frac{\lambda_{\max}(\bar{X}) + \alpha(\mu^* + \epsilon)}{\lambda_{\min}(\bar{X}) + \alpha} \leq \mu^* + 2\epsilon. \quad (14)$$

Recalling that  $X_\epsilon \succeq I$ , and  $\lambda_{\max}(\cdot)$  and  $\lambda_{\min}(\cdot)$  are convex and concave functions, respectively, we obtain that for any  $\alpha \geq \underline{\alpha}$ ,

$$\kappa(\bar{X} + \alpha X_\epsilon) = \frac{\lambda_{\max}((\bar{X} + \alpha X_\epsilon)/(1 + \alpha))}{\lambda_{\min}((\bar{X} + \alpha X_\epsilon)/(1 + \alpha))} \leq \frac{\lambda_{\max}(\bar{X}) + \alpha \lambda_{\max}(X_\epsilon)}{\lambda_{\min}(\bar{X}) + \alpha \lambda_{\min}(X_\epsilon)} \leq \frac{\lambda_{\max}(\bar{X}) + \alpha \lambda_{\max}(X_\epsilon)}{\lambda_{\min}(\bar{X}) + \alpha},$$

which together with (14) and the fact that  $\bar{X} + \alpha X_\epsilon \in \mathcal{S}_{++}^n \cap \Omega \forall \alpha \geq \underline{\alpha}$ , implies that

$$\kappa^* \leq \kappa(\bar{X} + \alpha X_\epsilon) \leq \mu^* + 2\epsilon \quad \forall \alpha \geq \underline{\alpha}. \quad (15)$$

Using (12), (15) and the arbitrariness of  $\epsilon$ , we conclude that  $\kappa^* \leq \mu^*$ , which together with the known result  $\mu^* \leq \kappa^*$  yields  $\kappa^* = \mu^*$ . Hence, statement (i) holds. Moreover, in view of (15) and the relation  $\kappa^* = \mu^*$ , we see that statement (iv) holds. Finally, recall from Theorem 3.1 (i) that  $\kappa^* = \lambda^*$ . Hence,  $\lambda^* = \mu^*$ . Using this relation and Theorem 3.1 (ii), we conclude that statement (ii) holds. ■

We next show that problem (10) is solvable, that is, its optimal value is achievable. Before proceeding, we provide some upper bounds on  $\epsilon$ -optimal solutions of (10) for some  $\epsilon > 0$ .

**Lemma 3.4** *Suppose that  $(t_\epsilon, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (10) for some  $\epsilon > 0$ . Define*

$$\underline{\lambda}^* = \inf \{ \lambda_{\max}(X) : X \in \mathcal{S}_+^n \cap \Omega \}. \quad (16)$$

Then

$$0 \leq t_\epsilon \leq (\mu^* + \epsilon)/\underline{\lambda}^*, \quad I \preceq X_\epsilon \preceq (\mu^* + \epsilon)I, \quad (17)$$

where  $\mu^*$  is the optimal value of (10).

*Proof.* By Assumptions A.1 and A.2, we observe that  $\underline{\lambda}^* \in (0, \infty)$ . Since  $(t_\epsilon, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (10), we know that  $\lambda_{\max}(X_\epsilon) \leq \mu^* + \epsilon$  and  $X_\epsilon \succeq I$ , which implies that the second relation of (17) holds. If  $t_\epsilon = 0$ , the first relation of (17) evidently holds. We now assume that  $t_\epsilon > 0$ . It follows from Theorem 3.3 (iii) that  $X_\epsilon/t_\epsilon \in \mathcal{S}_+^n \cap \Omega$ . This relation together with the definition of  $\underline{\lambda}^*$  implies that  $\lambda_{\max}(X_\epsilon)/t_\epsilon \geq \underline{\lambda}^*$ . Using this inequality and the relation  $\lambda_{\max}(X_\epsilon) \leq \mu^* + \epsilon$ , we see that the first relation of (17) holds. ■

We are now ready to show that problem (10) is solvable.

**Theorem 3.5** *Problem (10) is solvable, that is, its optimal value can be achieved at some feasible point.*

*Proof.* Given an arbitrary  $\epsilon > 0$ , define

$$\Pi := \{(t, X) : 0 \leq t \leq (\mu^* + \epsilon)/\underline{\lambda}^*, I \preceq X \preceq (\mu^* + \epsilon)I\},$$

where  $\mu^*$  is the optimal value of (10) and  $\underline{\lambda}^*$  is defined in (16). It follows from Lemma 3.4 that problem (10) is equivalent to

$$\inf \{ \lambda_{\max}(X) : (t, X) \in \Xi \cap \Pi \}, \quad (18)$$

where  $\Xi$  is defined in (11). We know that  $\Xi$  is closed. Hence,  $\Xi \cap \Pi$  is compact. In addition,  $\lambda_{\max}(\cdot)$  is continuous. It follows that problem (18) is solvable, which implies that problem (10) is solvable. ■

In view of Theorems 3.3 and 3.5, we see that problem (1) can be reformulated as a solvable convex programming problem (10). Moreover, an optimal solution of (10) can provide either an optimal or approximate solution of (1).

We now present three examples to illustrate how problem (1) can be solved as a convex programming problem. In the first two examples we choose  $\Omega$  to be the same sets as used in [8] for estimating the covariance matrix for the Markowitz portfolio selection model (see also [13]). In the third example we consider a positive semidefinite representable set  $\Omega$ . We show that for all these sets  $\Omega$ , problem (1) can be cast into a semidefinite programming (SDP) problem, which can be suitably solved by IP solvers (e.g., SeDuMi [16] and SDPT3 [17]) and first-order methods (see, for example, [6, 10] and Section 4).

**Corollary 3.6** *Let  $Q_1, \dots, Q_m \in \mathcal{S}^n$  be given. Assume that  $\Omega = \text{conv} \{Q_1, \dots, Q_m\}$ . Then problem (1) can be solved as the following SDP problem:*

$$\begin{aligned} \min_{s,y,X} \quad & s \\ \text{s.t.} \quad & \sum_{i=1}^m y_i Q_i - X = 0, \\ & y \in \mathfrak{R}_+^m, \quad I \preceq X \preceq sI. \end{aligned} \tag{19}$$

*Proof.* Since  $\Omega = \text{conv} \{Q_1, \dots, Q_m\}$ ,  $\Omega$  is a compact convex set and so  $\Omega_\infty = \{0\}$  by Proposition 2.2 (iii). Using this result, the definition of  $\Omega$  and Theorem 3.3, we see that problem (1) can be solved as the following SDP problem:

$$\begin{aligned} \min_{s,t,y,X} \quad & s \\ \text{s.t.} \quad & \sum_{i=1}^m y_i Q_i - X = 0, \\ & \sum_{i=1}^m y_i - t = 0, \\ & t \geq 0, \quad y \in \mathfrak{R}_+^m, \quad I \preceq X \preceq sI, \end{aligned}$$

which is equivalent to problem (19). Thus the conclusion holds. ■

**Corollary 3.7** *Let  $Q \in \mathcal{S}^n$  be given. Assume that  $\Omega = \{X \in \mathcal{S}^n : |X_{ij} - Q_{ij}| \leq \eta \forall ij\}$  for some  $\eta > 0$ . Then, problem (1) can be solved as the following SDP problem:*

$$\begin{aligned} \min_{s,t,X} \quad & s \\ \text{s.t.} \quad & (Q_{ij} - \eta)t \leq X_{ij} \leq (Q_{ij} + \eta)t \quad \forall ij \\ & t \geq 0, \quad I \preceq X \preceq sI. \end{aligned}$$

*Proof.* The conclusion follows from the definition of  $\Omega$  and Theorem 3.3. ■



**Corollary 3.8** Assume that  $\emptyset \neq \Omega$  is positive semidefinite representable, i.e., there exists  $C \in \mathcal{S}^m$  and linear operators  $\mathcal{A} : \mathcal{S}^n \rightarrow \mathcal{S}^m$  and  $\mathcal{B} : \mathfrak{R}^k \rightarrow \mathcal{S}^m$  such that

$$\Omega = \left\{ X \in \mathcal{S}^n : \mathcal{A}(X) + \mathcal{B}(u) + C \succeq 0 \text{ for some } u \in \mathfrak{R}^k \right\}. \quad (20)$$

Suppose that  $\mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^m$  is closed. Then, problem (1) can be solved as the following SDP problem:

$$\begin{aligned} \min_{s,t,u,X} \quad & s \\ \text{s.t.} \quad & \mathcal{A}(X) + \mathcal{B}(u) + tC \succeq 0, \\ & t \geq 0, \quad I \preceq X \preceq sI. \end{aligned} \quad (21)$$

*Proof.* First, notice that  $\Omega = \mathcal{A}^{-1}(\mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^m - C)$ . By the assumption that  $\mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^m$  is a closed convex cone, we see that  $\Omega$  is closed and convex. Moreover, it follows from the definition of recession cone that  $(\mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^m - C)_\infty = \mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^m$ . Using this relation and Proposition 2.2 (iv), we obtain that

$$\Omega_\infty = \mathcal{A}^{-1}((\mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^m - C)_\infty) = \mathcal{A}^{-1}(\mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^m). \quad (22)$$

Recall from Theorem 3.3 that problem (1) can be solved as (10), which together with (22) implies that the conclusion holds.  $\blacksquare$

*Remark.* If  $\mathfrak{R}(\mathcal{B}) \cap \mathcal{S}_+^n = \{0\}$  or  $\mathfrak{R}(\mathcal{B}) \cap \mathcal{S}_{++}^n \neq \emptyset$  holds, then  $\mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^n$  is closed (see [15]). Nevertheless, it generally may not be closed. For example, let  $\mathcal{B} : \mathfrak{R} \rightarrow \mathcal{S}^2$  be defined as:

$$\mathcal{B}(u) = \begin{bmatrix} u & 0 \\ 0 & 0 \end{bmatrix} \quad \forall u \in \mathfrak{R}.$$

Consider the sequence  $\{X_k\}$  defined as follows:

$$X_k = \begin{bmatrix} 0 & 1 \\ 1 & 1/k \end{bmatrix} = \begin{bmatrix} -k & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} k & 1 \\ 1 & 1/k \end{bmatrix} \quad \forall k \geq 1.$$

Then we have  $\{X_k\} \subseteq \mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^2$ , but

$$\lim_{k \rightarrow \infty} X_k = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \notin \mathfrak{R}(\mathcal{B}) + \mathcal{S}_+^2$$

since

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} - \begin{bmatrix} u & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} -u & 1 \\ 1 & 0 \end{bmatrix} \notin \mathcal{S}_+^2 \quad \forall u \in \mathfrak{R}.$$

Additionally, as pointed out by a referee, an SDP reformulation of the minimization of condition number is also discussed in [5, Exercise 4.43], where  $\Omega$  is an affine set given by  $\{Q_0 + \sum_{i=1}^m y_i Q_i : y \in \mathfrak{R}^m\}$  for some  $Q_0, \dots, Q_m \in \mathcal{S}^n$ .  $\blacksquare$

## 4 First-order method for finding minimum condition number

In this section, we propose a first-order method, the alternating direction method (ADM), for solving problem (1). In general, the ADM can be applied to solve problems of the following form:

$$\begin{aligned} \min_{x,y} \quad & f(x) + g(y) \\ \text{s.t.} \quad & \mathcal{A}x + \mathcal{B}y = b, \\ & x \in C_1, y \in C_2, \end{aligned} \tag{23}$$

where  $f$  and  $g$  are convex functions,  $\mathcal{A}$  and  $\mathcal{B}$  are linear operators, and  $C_1$  and  $C_2$  are closed convex sets. Each iteration of the ADM involves solving two subproblems successively, followed by an update of the multiplier. The method converges to an optimal solution of (23) under some mild assumptions (see, for example, [4, 7]).

### 4.1 Alternating direction method

In this subsection, we describe the implementation details of the ADM for solving problem (1). First, we notice from Theorem 3.3 that (1) can be reformulated as

$$\begin{aligned} \min_{X,t} \quad & \lambda_{\max}(X) \\ \text{s.t.} \quad & (t, X) \in \Xi, X \succeq I. \end{aligned} \tag{24}$$

Further, we can reformulate (24) as follows:

$$\begin{aligned} v^* := \min_{X,t,Y} \quad & \lambda_{\max}(X) \\ \text{s.t.} \quad & X - Y = 0, \\ & X \succeq I, (t, Y) \in \Xi. \end{aligned} \tag{25}$$

We then see that (25) is in the form of (23) with  $f = \lambda_{\max}(\cdot)$ ,  $g = 0$ ,  $\mathcal{A} = \mathcal{I}$ ,  $\mathcal{B} = \begin{pmatrix} 0 & -\mathcal{I} \end{pmatrix}$ ,  $b = 0$ ,  $C_1 = I + \mathcal{S}_+^n$  and  $C_2 = \Xi$ , where  $\mathcal{I}$  is the identity map. Thus, the ADM can be suitably applied to solve (25) (or, equivalently, (1)). To proceed, we introduce the following augmented Lagrangian function on  $\mathcal{S}^n \times \Re \times \mathcal{S}^n \times \mathcal{S}^n$ :

$$L_\beta(X, t, Y, \Gamma) = \lambda_{\max}(X) + \langle \Gamma, X - Y \rangle + \frac{\beta}{2} \|X - Y\|_F^2$$

for some  $\beta > 0$ .

We are now ready to present the algorithmic framework for the ADM when applied to solve (25) (or, equivalently, (1)).

#### Alternating direction method:

1. **Start:** Let  $(t^0, Y^0, \Gamma^0) \in \Re_+ \times \mathcal{S}^n \times \mathcal{S}^n$  and  $\beta > 0$  be given.

2. **For**  $k = 0, 1, \dots$

$$\begin{cases} \text{Compute } X^{k+1} \text{ by approximately solving } \min\{L_\beta(X, t^k, Y^k, \Gamma^k) : X \succeq I\}, \\ \text{Compute } (t^{k+1}, Y^{k+1}) \text{ by approximately solving } \min\{L_\beta(X^{k+1}, t, Y, \Gamma^k) : (t, Y) \in \Xi\}, \\ \Gamma^{k+1} = \Gamma^k + \beta(X^{k+1} - Y^{k+1}). \end{cases} \quad (26)$$

**End** (for)

Before discussing the convergence of the above method, we give the dual problem of (25) in the next proposition.

**Proposition 4.1** *The dual problem of (25) is given by*

$$\begin{aligned} \max_{\Gamma} \quad & 1 + \text{tr}(\Gamma) \\ \text{s.t.} \quad & e^T \max\{-\gamma, 0\} \leq 1, \Gamma \in \Omega^\ominus, \end{aligned} \quad (27)$$

where  $\gamma$  is the vector of eigenvalues of  $\Gamma$  and  $\Omega^\ominus$  is the negative polar of  $\Omega$ , i.e.,  $\Omega^\ominus = \{\Lambda \in \mathcal{S}^n : \langle \Lambda, Y \rangle \leq 0 \ \forall Y \in \Omega\}$ .

*Proof.* Recall that  $\Xi = \text{cl } \mathcal{K}$ . It then follows from Proposition 2.1 (iv) that  $\text{ri } \Xi = \text{ri}(\text{cl } \mathcal{K}) = \text{ri } \mathcal{K}$ . Using this relation and Proposition 2.1 (v), we obtain that

$$\text{ri} \{(X, t, Y) : X \succeq I, (t, Y) \in \Xi\} = \{(X, t, Y) : X \succ I, (t, Y) \in \text{ri } \mathcal{K}\}.$$

This equality together with (9) and the fact  $\mathcal{S}_{++}^n \cap \text{ri } \Omega \neq \emptyset$  implies that

$$\{(X, t, Y) : X - Y = 0\} \cap (\text{ri} \{(X, t, Y) : X \succeq I, (t, Y) \in \Xi\}) \neq \emptyset,$$

and hence, the Slater condition holds for (25). It then follows from this relation and [14, Corollary 28.2.2] that

$$\begin{aligned} v^* &= \min_{X, Y} \{\lambda_{\max}(X) : X = Y, X \succeq I, (t, Y) \in \Xi\} = \min_{X \succeq I, (t, Y) \in \Xi} \max_{\Gamma} \{\lambda_{\max}(X) + \langle \Gamma, X - Y \rangle\} \\ &= \max_{\Gamma} \min_{X \succeq I, (t, Y) \in \Xi} \{\lambda_{\max}(X) + \langle \Gamma, X - Y \rangle\}. \end{aligned}$$

Using the fact that  $\lambda_{\max}(X) = \max\{\langle P, X \rangle : \text{tr}(P) = 1, P \succeq 0\}$ , we see further that

$$\begin{aligned} v^* &= \max_{\Gamma} \min_{X \succeq I, (t, Y) \in \Xi} \max_{\text{tr}(P)=1, P \succeq 0} \{\langle P, X \rangle + \langle \Gamma, X - Y \rangle\} \\ &= \max_{\text{tr}(P)=1, P \succeq 0, \Gamma} \min_{X \succeq I, (t, Y) \in \Xi} \{\langle P + \Gamma, X \rangle + \langle -\Gamma, Y \rangle\} \\ &= \max_{P, \Gamma} \{1 + \text{tr}(\Gamma) : \text{tr}(P) = 1, P \succeq 0, P + \Gamma \succeq 0, \Gamma \in (P_Y(\Xi))^\ominus\}, \end{aligned} \quad (28)$$

where the second equality follows from [14, Corollary 37.3.2] and  $P_Y(\Xi)$  is the projection of  $\Xi$  onto the  $Y$ -coordinate. Furthermore, we obtain from (11) and Proposition 2.2 (i) that  $P_Y(\Xi) = \text{cone}(\Omega) \cup \Omega_\infty = \text{cl } \text{cone}(\Omega)$ . Using this result and the definition of negative polar, we have

$$(P_Y(\Xi))^\ominus = (\text{cl}(\text{cone}(\Omega)))^\ominus = \Omega^\ominus. \quad (29)$$

Finally, we observe that  $-\Gamma \preceq P$  holds for some  $P$  satisfying  $\text{tr}(P) = 1$  and  $P \succeq 0$  if and only if  $e^T \max\{-\gamma, 0\} \leq 1$ , where  $\gamma$  is the vector of eigenvalues of  $\Gamma$ . The conclusion of this proposition immediately follows from this observation, (28) and (29).  $\blacksquare$

We are now ready to state a convergence result for the above ADM, which asserts global convergence provided that the subproblems are solved sufficiently accurately. Its proof can be found in [7, Theorem 8].

**Proposition 4.2** *Let  $\beta > 0$  and  $\{\nu_k\}$  be a sequence of nonnegative numbers with  $\sum \nu_k < \infty$ . Let  $\{(X^k, t^k, Y^k, \Gamma^k)\}$  be generated as in (26) with  $\{X^k\}$  and  $\{(t^k, Y^k)\}$  satisfying*

$$\left\| X^k - \underset{X \succeq I}{\text{argmin}} L_\beta(X, t^{k-1}, Y^{k-1}, \Gamma^{k-1}) \right\|_F \leq \nu_k,$$

$$\inf \left\{ \left\| Y^k - Y \right\|_F : (t, Y) \in \underset{(t, Y) \in \Xi}{\text{Argmin}} L_\beta(X^k, t, Y, \Gamma^{k-1}) \right\} \leq \nu_k$$

for all  $k$ . Then  $\{(X^k, Y^k, \Gamma^k)\}$  is convergent. Furthermore, any accumulation point of  $\{(X^k, t^k)\}$  solves (24) and the limit of  $\{\Gamma^k\}$  solves (27).

We next discuss how the two subproblems in (26) can be solved efficiently. We start by considering the first subproblem, namely,  $\min\{L_\beta(X, t^k, Y^k, \Gamma^k) : X \succeq I\}$ . Notice that this subproblem can be formulated as follows:

$$\min_X \left\{ \frac{1}{\beta} \lambda_{\max}(X) + \frac{1}{2} \left\| X - \left( Y^k - \frac{\Gamma^k}{\beta} \right) \right\|_F^2 : X \succeq I \right\}. \quad (30)$$

Let  $U^T \text{diag}(\xi^k) U$  be an eigenvalue decomposition of  $Y^k - \frac{\Gamma^k}{\beta} - I$ , where  $U \in \mathfrak{R}^{n \times n}$  is an orthogonal matrix, and  $\text{diag}(\xi^k)$  is a diagonal matrix whose diagonal consists of the vector  $\xi^k$ . Since  $\lambda_{\max}(\cdot)$ ,  $\|\cdot\|_F$  and  $\{X : X \succeq I\}$  are unitary similarity invariant, it follows from [12, Proposition 2.7] that the solution  $X^*$  of (30) is given by

$$X^* = U^T \text{diag}(x^*) U + I,$$

where  $x^* \in \mathfrak{R}^n$  is the optimal solution of

$$v^k := \min_x \left\{ f^k(x) := \frac{1}{\beta} \max_i x_i + \frac{1}{2} \|x - \xi^k\|^2 : x \geq 0 \right\}. \quad (31)$$

Problem (31) has a nonsmooth objective function. Its dual problem, however, is smooth as shown in the following proposition.

**Proposition 4.3** *The dual problem of (31) is given by*

$$\begin{aligned} \max_w \quad & d^k(w) := \frac{1}{2} \|\xi^k\|^2 - \frac{1}{2} \|(\xi^k - w)_+\|^2 \\ \text{s.t.} \quad & e^T w = \frac{1}{\beta}, w \geq 0. \end{aligned} \quad (32)$$

Furthermore, if  $w^*$  solves (32), then  $(\xi^k - w^*)_+$  solves (31).

*Proof.* Note that we have

$$\begin{aligned} v^k &= \min_{x \geq 0} \left\{ \frac{1}{\beta} \max_i x_i + \frac{1}{2} \|x - \xi^k\|^2 \right\} = \min_{x \geq 0} \max_{e^T w = \frac{1}{\beta}, w \geq 0} \left\{ w^T x + \frac{1}{2} \|x - \xi^k\|^2 \right\} \\ &= \max_{e^T w = \frac{1}{\beta}, w \geq 0} \min_{x \geq 0} \left\{ w^T x + \frac{1}{2} \|x - \xi^k\|^2 \right\} = \max_{e^T w = \frac{1}{\beta}, w \geq 0} \left\{ \frac{1}{2} \|\xi^k\|^2 - \frac{1}{2} \|(\xi^k - w)_+\|^2 \right\}, \end{aligned}$$

where the third equality follows from [14, Corollary 37.3.2], and the inner minimization in this equality is achieved at  $x = (\xi^k - w)_+$ . The conclusion of the proposition immediately follows.  $\blacksquare$

In view of Proposition 4.3, the solution of (30) can be found by solving (32). Since the objective function of (32) is smooth and the projection onto simplices can be computed efficiently (see [18]), the spectral projected gradient (SPG) method (see, for example, [3, 11]) can be suitably applied to solve (32).

We next discuss how to solve the second subproblem of (26), namely,  $\min\{L_\beta(X^{k+1}, t, Y, \Gamma^k) : (t, Y) \in \Xi\}$ . Notice that this subproblem can be formulated as

$$\min_{(t, Y) \in \Xi} \frac{1}{2} \left\| Y - \left( X^{k+1} + \frac{\Gamma^k}{\beta} \right) \right\|_F^2. \quad (33)$$

Observe from the definition of  $\Xi$  that  $\{Y : (t, Y) \in \Xi\} = \text{cone}(\Omega) \cup \Omega_\infty$ . Since  $0 \notin \Omega$ , it then follows from Proposition 2.2 (i) that  $\{Y : (t, Y) \in \Xi\} = \text{cl cone}(\Omega)$ . Hence, the optimal solution  $Y^*$  of (33) is the projection of  $X^{k+1} + \frac{\Gamma^k}{\beta}$  onto  $\text{cl cone}(\Omega)$ , and  $t^*$  is any nonnegative number such that  $(t^*, Y^*) \in \Xi$ . For a general  $\Omega$ , problem (33) does not have a closed form solution. In the remainder of this subsection, we will discuss how (33) can be solved efficiently for  $\Omega$  as specified in Corollaries 3.6 and 3.7, and also identify  $\Omega^\ominus$  explicitly that is used in (27). Since  $\Omega$  is compact in these cases, it follows from (11) and Proposition 2.2 (iii) that

$$\Xi = \{t(1, X) : t \geq 0, X \in \Omega\}. \quad (34)$$

We first consider the case where  $\Omega$  is the convex hull of some symmetric matrices, i.e.,

$$\Omega = \text{conv}\{Q_1, \dots, Q_m\} := \left\{ \sum_{i=1}^m y_i Q_i : \sum_{i=1}^m y_i = 1, y_i \in [0, 1] \right\} \quad (35)$$

for some  $Q_1, \dots, Q_m \in \mathcal{S}^n$ . For such  $\Omega$ , problem (33) reduces to

$$\min_{y \geq 0} H^k(y) := \frac{1}{2} \left\| \sum_{i=1}^m y_i Q_i - \left( X^{k+1} + \frac{\Gamma^k}{\beta} \right) \right\|_F^2. \quad (36)$$

Suppose that  $y^*$  is an optimal solution of (36). Then  $(t^*, Y^*) = (\sum_{i=1}^m y_i^*, \sum_{i=1}^m y_i^* Q_i)$  is an optimal solution of (33). Since the objective function of (36) is smooth and projection onto the nonnegative orthant is cheap, problem (36) can be suitably solved by the SPG method. In addition, the negative polar of  $\Omega$  is given by

$$\Omega^\ominus = \{\Gamma \in \mathcal{S}^n : \langle \Gamma, Y \rangle \leq 0 \forall Y \in \Omega\} = \{\Gamma \in \mathcal{S}^n : \text{tr}(Q_i \Gamma) \leq 0 \forall i = 1, \dots, m\}.$$

Finally, we consider the case where  $\Omega$  is a neighborhood of some  $Q \in \mathcal{S}^n$ , i.e.,

$$\Omega = \{Y \in \mathcal{S}^n : |Y_{ij} - Q_{ij}| \leq \eta \forall ij\} = \{Y \in \mathcal{S}^n : Q_{ij} - \eta \leq Y_{ij} \leq Q_{ij} + \eta \forall ij\} \quad (37)$$

for some  $\eta > 0$ . From (34), we see that problem (33) reduces to

$$\begin{aligned} & \min_{t \geq 0} \min_{Y \in \Omega} \frac{1}{2} \left\| Y - \left( X^{k+1} + \frac{\Gamma^k}{\beta} \right) \right\|_F^2 \\ &= \min_{t \geq 0} \frac{1}{2} \left\| \min \left\{ \max \left\{ X^{k+1} + \frac{\Gamma^k}{\beta}, (Q - \eta E)t \right\}, (Q + \eta E)t \right\} - \left( X^{k+1} + \frac{\Gamma^k}{\beta} \right) \right\|_F^2. \end{aligned} \quad (38)$$

We observe that if  $t^*$  is an optimal solution of problem (38), then  $(t^*, Y^*)$  is an optimal solution of (33), where

$$Y^* = \min \left\{ \max \left\{ X^{k+1} + \frac{\Gamma^k}{\beta}, (Q - \eta E)t^* \right\}, (Q + \eta E)t^* \right\}.$$

Notice that the objective function of (38) is smooth. For the similar reason as mentioned above, problem (38) can be suitably solved by the SPG method. In addition, the negative polar  $\Omega^\ominus$  is given by

$$\begin{aligned} \Omega^\ominus &= \{\Gamma \in \mathcal{S}^n : \langle \Gamma, Y \rangle \leq 0 \forall Y \in \Omega\} = \left\{ \Gamma \in \mathcal{S}^n : \max_{Y \in \Omega} \text{tr}(\Gamma Y) \leq 0 \right\} \\ &= \left\{ \Gamma \in \mathcal{S}^n : \sum_{\Gamma_{ij} > 0} \Gamma_{ij}(Q_{ij} + \eta) + \sum_{\Gamma_{ij} < 0} \Gamma_{ij}(Q_{ij} - \eta) \leq 0 \right\}. \end{aligned}$$

## 4.2 Computational results

In this section, we conduct numerical experiments to test the performance of our approach on solving (24) with different  $\Omega$ . In particular, we compare our method with the standard interior point (IP) solver SeDuMi1.1R3 [16]. Our codes for ADM are written in Matlab while SeDuMi is coded in C with a Matlab interface. All experiments are performed in Matlab Version 7.8 on a Dell POWEREDGE 1950 with Debian 5.0.6 (Linux).

### 4.2.1 $\Omega$ given by (35)

We consider numerical experiments with  $\Omega$  given by (35). In particular, we choose  $\Omega$  to be the convex hull of some positive definite matrices. For the ADM, we set  $t^0 = 0$ ,  $Y^0 = 0$ ,  $\Gamma^0 = -I/n$ , and terminate the method once

$$\begin{aligned} \frac{|\lambda_{\max}(Y^k) - 1 - \text{tr}(\Gamma^k)|}{\max\{1, |\lambda_{\max}(Y^k)|\}} &< 10 \cdot \text{tol}, & \frac{1 - \lambda_{\min}(Y^k)}{\max\{1, \|Y^k\|\}} &< \text{tol}, \\ \frac{\max\{e^T \max\{-\gamma^k, 0\} - 1, \max_{1 \leq i \leq m} \{\text{tr}(Q_i \Gamma^k)\}\}}{\max\{1, \|\Gamma^k\|_F\}} &< \text{tol} \end{aligned}$$

Table 1: Computational results for solving (24) with  $\Omega$  given by (35)

$n$	$m$	$\min_{1 \leq i \leq m} \kappa(Q_i)$	cond		cpu	
			SeDuMi	ADM	SeDuMi	ADM
50	80	1.828e+03	1.44	1.44	1.85	1.64
50	100	1.789e+03	1.38	1.38	2.34	1.95
50	120	1.730e+03	1.33	1.33	3.08	2.58
60	80	2.774e+03	1.46	1.46	2.60	3.21
60	100	2.270e+03	1.39	1.39	3.43	3.59
60	120	2.503e+03	1.35	1.35	4.33	3.75

for some  $tol > 0$ , where  $\gamma^k$  is the vector of eigenvalues of  $\Gamma^k$ . In addition, we apply the SPG method to find approximate solutions  $w^k$  and  $y^k$  to (32) and (36), respectively, and terminate the method once

$$\frac{|f^k((\xi^k - w^k)_+) - d^k(u^l)|}{\max\{1, |d^k(w^k)|\}} < \min\{1e - 6, tol/100\},$$

$$\frac{|\max\{y^k - \nabla H^k(y^k), 0\} - y^k|}{\max\{1, |H^k(y^k)|\}} < \min\{1e - 6, tol/100\}.$$

Suppose that  $(t^k, X^k, Y^k)$  is the approximate solution obtained by the ADM. Clearly,  $Y^k/t^k \in \mathcal{S}_+^n \cap \Omega$ . We then compute an approximate minimum condition number as  $\kappa(Y^k/t^k)$ . For SeDuMi, we use the default tolerance and compute an approximate minimum condition number similarly.

In our experiments, for each  $n = 50, 60$  and each  $m = 80, 100, 120$ , we generate 10 samples of  $n \times n$  matrices  $B_i$ ,  $i = 1, \dots, m$ , with i.i.d. Gaussian entries. We then set  $Q_i = B_i B_i^T$  for each  $i = 1, \dots, m$  and solve the corresponding problem (24) with  $\Omega$  defined as in (35). We set  $\beta = 1$  and  $tol = 1e - 4$  for the ADM. The results of this experiment are reported in Table 1. In particular, we report the CPU time (cpu) in seconds and the approximate minimum condition number (cond) computed as above for ADM and SeDuMi, averaged over 10 instances. We also report  $\min_{1 \leq i \leq m} \kappa(Q_i)$  for each  $(n, m)$ , averaged over the 10 instances. We see from Table 1 that our method is comparable with SeDuMi in terms of CPU time, and gives the same approximate minimum condition number.

#### 4.2.2 $\Omega$ given by (37)

We consider numerical experiments with  $\Omega$  given by (37). For the ADM, we set  $t^0 = 1$ ,  $Y^0 = I$ ,  $\Gamma^0 = -I/n$ , and terminate the method once  $Y^k \succ 0$ , and

$$\frac{|\lambda_{\max}(Y^k) - 1 - \text{tr}(\Gamma^k)|}{\max\{1, |\lambda_{\max}(Y^k)|\}} < 10 \cdot tol, \quad \frac{1 - \lambda_{\min}(Y^k)}{\max\{1, \|Y^k\|\}} < tol,$$

$$\frac{\max\left\{e^T \max\{-\gamma^k, 0\} - 1, \sum_{\Gamma_{ij}^k < 0} \Gamma_{ij}^k (Q_{ij} - \eta) + \sum_{\Gamma_{ij}^k > 0} \Gamma_{ij}^k (Q_{ij} + \eta)\right\}}{\max\{1, \|\Gamma^k\|_F\}} < tol$$

for some  $tol > 0$ . In addition, we apply the SPG methods similarly as in the previous subsection to approximately solve (32) and (38). Suppose that  $(t^k, X^k, Y^k)$  is the approximate solution obtained by the ADM. Clearly,  $Y^k/t^k \in \mathcal{S}_+^n \cap \Omega$ . We then compute an approximate minimum condition number as

Table 2: Computational results for solving (24) with  $\Omega$  given by (37)

$n$	$\eta$	$\kappa(Q)$	cond		cpu	
			SeDuMi	ADM	SeDuMi	ADM
50	0.5	4.456e+06	28.4	28.5	23.6	6.3
50	1	2.779e+08	15.9	16.0	26.1	5.9
60	0.5	2.936e+05	30.2	30.2	98.1	8.9
60	1	7.607e+04	16.3	16.3	105.5	4.4
70	0.5	1.071e+07	31.8	31.8	297.7	14.6
70	1	1.486e+06	18.0	18.0	307.4	6.7

$\kappa(Y^k/t^k)$ . For SeDuMi, we use the default tolerance and compute an approximate minimum condition number similarly.

In our experiments, for each  $n = 50, 60, 70$  and each  $\eta = 0.5, 1$ , we generate 10 samples of  $n \times n$  matrix  $A$  with i.i.d. Gaussian entries and set  $Q = AA^T$ . We then solve the corresponding problem (24) with  $\Omega$  defined as in (37). We set  $\beta = 0.1$  and  $tol = 1e - 4$  for the ADM. The results of this experiment are reported in Table 2. In particular, we report the CPU time (cpu) in seconds and the approximate minimum condition number (cond) computed as above for both methods, averaged over 10 instances. We also report  $\kappa(Q)$  for each  $n$ , averaged over the 10 instances. We see from Table 2 that our method is significantly faster than SeDuMi. In addition, our method sometimes yields a slightly larger condition number than SeDuMi, but the difference is negligible relative to  $\kappa(Q)$ .

## 5 Finding optimal preconditioners

Preconditioning techniques have been widely used in solving large-scale linear systems arising in many applications. In recent years, various approaches have been proposed for finding practical preconditioners such as incomplete factorization methods, sparse approximate inverses and more recently some other variants based on the multilevel paradigm. We refer readers to [2] and the references therein for a comprehensive survey. Instead of proposing another efficient method for finding good preconditioners, we study the following question: what could be the best one in a set of infinitely many preconditioners? In particular, we consider the optimal preconditioner finding problem, that is, problem (4).

We first show that problem (4) can be solved as the following problem:

$$\inf \{ \lambda_{\max}(X^T C^T C X) : (t, X) \in \Xi, X^T C^T C X \succeq I \}, \quad (39)$$

where  $\Xi$  is defined as in (11).

**Proposition 5.1** *The following statements hold:*

- i) problem (39) has the same optimal value as (4);*
- ii) for any  $\epsilon \geq 0$ , if  $X_\epsilon$  is an  $\epsilon$ -optimal solution of (4), then  $(1/\delta_\epsilon, X_\epsilon/\delta_\epsilon)$  is an  $\epsilon$ -optimal solution of (39), where  $\delta_\epsilon = \sqrt{\lambda_{\min}(X_\epsilon^T C^T C X_\epsilon)}$ ;*
- iii) for any  $\epsilon > 0$ , if  $(t_\epsilon, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (39) for some  $t_\epsilon > 0$ , then  $X_\epsilon/t_\epsilon$  is an  $\epsilon$ -optimal solution of (4);*



iv) for any  $\epsilon > 0$ , if  $(0, X_\epsilon)$  is an  $\epsilon$ -optimal solution of (39), then  $\bar{X} + \alpha X_\epsilon$  is an  $2\epsilon$ -optimal solution of (4), provided that  $\alpha$  is sufficiently large, where  $\bar{X}$  is an arbitrary matrix in  $\Omega$ .

*Proof.* Given any  $X, Y \in \mathfrak{R}^{n \times n}$ , we can see that

$$\begin{aligned} \lim_{\alpha \rightarrow \infty} \kappa((X + \alpha Y)^T C^T C (X + \alpha Y)) &= \lim_{\alpha \rightarrow \infty} \frac{\lambda_{\max}((X + \alpha Y)^T C^T C (X + \alpha Y))}{\lambda_{\min}((X + \alpha Y)^T C^T C (X + \alpha Y))} \\ &= \lim_{\alpha \rightarrow \infty} \frac{\lambda_{\max}((X/\alpha + Y)^T C^T C (X/\alpha + Y))}{\lambda_{\min}((X/\alpha + Y)^T C^T C (X/\alpha + Y))} \\ &= \lim_{\alpha \rightarrow \infty} \frac{\lambda_{\max}(Y^T C^T C Y)}{\lambda_{\min}(Y^T C^T C Y)} = \kappa(Y^T C^T C Y), \end{aligned} \quad (40)$$

where the third inequality is due to the continuity of eigenvalues. The conclusion of this proposition then follows from (40) and similar arguments as used in the proof of Theorem 3.3.  $\blacksquare$

Notice that the set defined by the constraint  $X^T C^T C X \succeq I$  is typically nonconvex. Thus problem (39) is in general nonconvex. We next consider the special case where  $X$  is further restricted to be in  $\mathcal{S}_+^n$ , i.e.,

$$v_{\text{opt}} = \inf \{ \lambda_{\max}(X^T C^T C X) : (t, X) \in \Xi, X^T C^T C X \succeq I, X \in \mathcal{S}_+^n \}, \quad (41)$$

and propose a convex relaxation to this problem, which turns out to be exact when finding optimal diagonal preconditioners.

**Theorem 5.2** Consider the following convex programming problem:

$$\begin{aligned} v_{\text{relax}} &= \min_{s, t, Y, X} s \\ \text{s.t.} \quad &\begin{bmatrix} I & CX \\ X^T C^T & sI \end{bmatrix} \succeq 0, \\ &\begin{bmatrix} I & Y \\ Y & C^T C \end{bmatrix} \succeq 0, \\ &\begin{bmatrix} X & I \\ I & Y \end{bmatrix} \succeq 0, \\ &(t, X) \in \Xi, X \in \mathcal{S}_+^n, Y \in \mathcal{S}_+^n. \end{aligned} \quad (42)$$

Suppose  $(s^*, t^*, Y^*, X^*)$  is an optimal solution of (42). Then we have

$$\frac{v_{\text{relax}}}{\delta^*} \geq v_{\text{opt}} \geq v_{\text{relax}}, \quad (43)$$

where  $\delta^* = \lambda_{\min}(X^* Y^{*2} X^*)$ . Moreover, if we further require  $X, Y \in \mathcal{D}_+^n$  in (41) and (42) (that is, when the optimal diagonal preconditioner is sought), then the convex problem (42) is equivalent to (41).

*Proof.* We first observe that

$$s \geq \lambda_{\max}(X^T C^T C X) \Leftrightarrow sI - X^T C^T C X \succeq 0 \Leftrightarrow \begin{bmatrix} I & CX \\ X^T C^T & sI \end{bmatrix} \succeq 0.$$

In addition, for any  $X \in \mathcal{S}_+^n$ ,

$$\begin{aligned} X^T C^T C X \succeq I &\Leftrightarrow C^T C \succeq X^{-2} \Leftrightarrow C^T C \succeq Y^2, Y^2 \succeq X^{-2} \text{ for some } Y \in \mathcal{S}_+^n, \\ &\Rightarrow C^T C \succeq Y^2, Y \succeq X^{-1} \text{ for some } Y \in \mathcal{S}_+^n, \end{aligned} \quad (44)$$

$$\Leftrightarrow \begin{bmatrix} I & Y \\ Y & C^T C \end{bmatrix} \succeq 0, \begin{bmatrix} X & I \\ I & Y \end{bmatrix} \succeq 0 \quad (45)$$

where (44) follows from [9, Problem 6.6.17]. Using the above observations, we see that  $v_{\text{opt}} \geq v_{\text{relax}}$ . Furthermore, let  $(s^*, t^*, Y^*, X^*)$  be an optimal solution of (42). It follows from (44) and (45) that  $X^*, Y^* \in \mathcal{S}_{++}^n$  and

$$X^{*T} C^T C X^* \succeq X^* Y^{*2} X^* \succeq \delta^* I \succ 0. \quad (46)$$

Using (46) and the fact that  $\Xi$  is a cone, we further observe that  $(t^*, X^*)/\sqrt{\delta^*}$  is feasible for (41). We thus have

$$\frac{v_{\text{relax}}}{\delta^*} = \lambda_{\max} \left( \frac{X^*}{\sqrt{\delta^*}} C^T C \frac{X^*}{\sqrt{\delta^*}} \right) \geq v_{\text{opt}}$$

and hence (43) holds. Finally, if  $X$  and  $Y$  are further restricted to be in  $\mathcal{D}_+^n$ , then the implication in (44) becomes an equivalence, and thus (41) is equivalent to the convex programming problem (42). ■

The next proposition provides a lower bound on the value of  $\delta^*$ .

**Proposition 5.3** *For any optimal solution  $(s^*, t^*, Y^*, X^*)$  of (42), it holds that*

$$\delta^* = \lambda_{\min}(X^* Y^{*2} X^*) \geq \frac{1}{\sqrt{\kappa(C^T C) v_{\text{relax}}}},$$

where  $v_{\text{relax}}$  is the optimal value of (42).

*Proof.* First, we know from the proof of Theorem 5.2 that  $X^* \succeq Y^{*-1}$ , which implies that  $X^{*\frac{1}{2}} Y^* X^{*\frac{1}{2}} \succeq I$ . We thus have

$$\begin{aligned} \delta^* &= \lambda_{\min}(X^* Y^{*2} X^*) \geq \lambda_{\min}(X^* Y^* X^*) \lambda_{\min}(Y^*) \\ &= \lambda_{\min}(X^{*\frac{1}{2}} X^{*\frac{1}{2}} Y^* X^{*\frac{1}{2}} X^{*\frac{1}{2}}) \lambda_{\min}(Y^*) \\ &\geq \lambda_{\min}(X^*) \lambda_{\min}(Y^*). \end{aligned} \quad (47)$$

Next, from the second constraint in (42), we know that  $C^T C \succeq Y^{*2}$  and hence

$$\lambda_{\max}(C^T C) \geq (\lambda_{\max}(Y^*))^2. \quad (48)$$

Furthermore, recall from the proof of Theorem 5.2 that

$$v_{\text{relax}} = \lambda_{\max}(X^{*T} C^T C X^*) \geq \lambda_{\min}(C^T C) (\lambda_{\max}(X^*))^2. \quad (49)$$

Finally, the third constraint in (42) implies that  $X^* \succeq Y^{*-1}$  and  $Y^* \succeq X^{*-1}$ . Thus we have

$$\lambda_{\min}(X^*) \geq \lambda_{\min}(Y^{*-1}) = \frac{1}{\lambda_{\max}(Y^*)}, \quad \lambda_{\min}(Y^*) \geq \lambda_{\min}(X^{*-1}) = \frac{1}{\lambda_{\max}(X^*)}. \quad (50)$$

Combining (47)-(50), we obtain

$$\delta^* \geq \lambda_{\min}(X^*)\lambda_{\min}(Y^*) \geq \frac{1}{\lambda_{\max}(Y^*)\lambda_{\max}(X^*)} \geq \frac{1}{\sqrt{\kappa(C^T C)v_{\text{relax}}}}.$$

This completes the proof.  $\blacksquare$

We next present an example to illustrate how the problem of finding a diagonal preconditioner can be solved as a convex programming problem. In particular, we choose  $\Omega$  to be a box. One can see from Theorem 5.2 that such a problem can be cast into a cone programming problem, which can be suitably solved by interior point solvers (e.g., SeDuMi [16] and SDPT3 [17]) and first-order methods (see, for example, [10]).

**Corollary 5.4** *Let  $d \in \mathfrak{R}_{++}^n$  be given. Assume that  $\Omega = \{X \in \mathcal{D}_+^n : |X_{ii} - d_i| \leq \eta \forall i\}$  for some  $\eta > 0$ . Then, problem (4) can be solved as the following cone programming problem:*

$$\begin{aligned} \min_{s,t,Y,X} \quad & s \\ \text{s.t.} \quad & \begin{bmatrix} I & CX \\ X^T C^T & sI \end{bmatrix} \succeq 0, \\ & \begin{bmatrix} I & Y \\ Y & C^T C \end{bmatrix} \succeq 0, \\ & \begin{pmatrix} (Y_{ii} + X_{ii})/2 \\ (Y_{ii} - X_{ii})/2 \\ 1 \end{pmatrix} \in \mathcal{L}^3, \quad i = 1, \dots, n, \\ & (d_i - \eta)t \leq X_{ii} \leq (d_i + \eta)t, \quad i = 1, \dots, n, \\ & t \geq 0, X \in \mathcal{D}_+^n, Y \in \mathcal{D}_+^n. \end{aligned} \tag{51}$$

*Proof.* Notice that when  $X, Y \in \mathcal{D}_+^n$ , we have

$$\begin{bmatrix} X & I \\ I & Y \end{bmatrix} \succeq 0 \Leftrightarrow Y \succeq X^{-1} \Leftrightarrow Y_{ii}X_{ii} \geq 1 \forall i = 1, \dots, n \Leftrightarrow \begin{pmatrix} (Y_{ii} + X_{ii})/2 \\ (Y_{ii} - X_{ii})/2 \\ 1 \end{pmatrix} \in \mathcal{L}^3, \quad i = 1, \dots, n.$$

The conclusion of this corollary follows from this observation and Theorem 5.2.  $\blacksquare$

Before ending this section, we perform numerical experiments to compare the quality of the diagonal preconditioner obtained by solving (51) and the Jacobi diagonal preconditioner, i.e., setting  $X_{ii} = \frac{1}{\sqrt{(C^T C)_{ii}}}$ . For each  $n = 40, 50, 60$  and  $\eta = 0.05, 0.1$ , we generate 10 samples of  $n \times n$  matrix  $A$  with i.i.d. Gaussian entries. We then set  $Q = AA^T$ . We take  $C^T C$  to be the Cholesky decomposition of  $Q$  and set  $d_i = \frac{1}{\sqrt{(C^T C)_{ii}}}$  for  $i = 1, \dots, n$ . We then solve the corresponding problem (51) using the IP solver SeDuMi with the default tolerance. The approximate optimal condition number after preconditioning is computed as

$$\kappa_{\text{opt}} = \kappa \left( \frac{X^*}{t^*} C^T C \frac{X^*}{t^*} \right),$$

Table 3: Computational results for solving (51)

$n$	$\eta$	$\kappa(C^T C)$	$\kappa_J$	$\kappa_{\text{opt}}$
40	0.05	7.102e+04	6.871e+04	4.905e+04
40	0.10	8.810e+04	7.657e+04	5.848e+04
50	0.05	6.158e+05	5.672e+05	4.240e+05
50	0.10	7.120e+06	6.573e+06	4.954e+06
60	0.05	4.658e+06	4.344e+06	3.195e+06
60	0.10	4.884e+04	4.478e+04	3.196e+04

where  $(s^*, t^*, X^*, Y^*)$  is the approximate optimal solution given by SeDuMi. We report in Table 3 the original condition number  $\kappa(C^T C)$ , the condition number upon applying Jacobi diagonal preconditioner (denoted by  $\kappa_J$ ) and the approximate optimal condition number  $\kappa_{\text{opt}}$ , averaged over the 10 instances. We see that  $\kappa_{\text{opt}}$  is usually smaller than  $J$  by at least 20%.

## 6 Concluding remarks

In this paper we considered minimizing the spectral condition number of a positive semidefinite matrix over a nonempty closed convex set  $\Omega$ . We showed that it can be solved as a convex programming problem, and moreover the optimal value of the latter problem is achievable. As a consequence, when  $\Omega$  is positive semidefinite representable, it can be cast into an SDP problem. We also considered a closely related problem, that is, finding an optimal preconditioner for a positive definite matrix. We proposed a convex relaxation for finding positive definite preconditioners. This relaxation turns out to be exact when finding diagonal preconditioners.

The results of this paper can be extended to the problem:

$$\inf \left\{ \frac{\sum_{i=1}^k \lambda_i(X)}{\sum_{j=0}^l \lambda_{n-j}(X)} : X \in \mathcal{S}_+^n \cap \Omega \right\},$$

where  $1 \leq k \leq n$ ,  $0 \leq l \leq n - 1$ , and  $\lambda_i(X)$  denotes the  $i$ th largest eigenvalue of  $X$  for  $i = 1, \dots, n$ .

## Acknowledgement

The first author would like to thank Professor Jane Ye for bringing his attention to the topic of this paper and also for showing the results of the paper [13]. The second author would like to thank Dr Guoyin Li for discussion of more efficient implementations of SeDuMi. We would also like to thank the anonymous referees for their comments and for pointing out the reference [5].

## References

- [1] A. Auslender and M. Teboulle. *Asymptotic Cones and Functions in Optimization and Variational Inequalities*. Springer, 2003.
- [2] M. Benzi. Preconditioning techniques for large linear systems: a survey. *J. Comput. Phys.*, 182:418–477, 2002.

- [3] E. G. Birgin, J. M. Martínez and M. Raydan. Nonmonotone spectral projected gradient methods on convex sets. *SIAM J. Optim.*, 10:1196–1211, 2000.
- [4] D. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Prentice Hall, 1989.
- [5] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [6] S. Burer and R. D. C. Monteiro. A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Math. Program.*, 95:329–357, 2003.
- [7] J. Eckstein and D. Bertsekas. On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators. *Math. Program.*, 55:293–318, 1992.
- [8] V. Guigues. Inférence Statistique pour l’Optimisation Stochastique. Ph.D. thesis, Université Joseph Fourier, Grenoble, France, 2005.
- [9] R. A. Horn and C. R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press, 2009.
- [10] G. Lan, Z. Lu and R. D. C. Monteiro. Primal-dual first-order methods with  $O(1/\epsilon)$  iteration iteration-complexity for cone programming. *Math. Program.*, 126:1–29, 2011.
- [11] Z. Lu and Y. Zhang. An augmented Lagrangian approach for sparse principal component analysis. Technical Report, Department of Mathematics, Simon Fraser University, Burnaby, BC, V5A 1S6, Canada, July 2009.
- [12] Z. Lu and Y. Zhang. Penalty decomposition methods for rank minimization. Technical Report, Department of Mathematics, Simon Fraser University, Burnaby, BC, V5A 1S6, Canada, September 2010.
- [13] P. Maréchal and J. J. Ye. Optimizing condition numbers. *SIAM J. Optim.*, 20:935–947, 2009.
- [14] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- [15] G. Pataki. On the closedness of the linear image of a closed convex cone. *Math. Oper. Res.*, 32:395–412, 2007.
- [16] J. F. Sturm. Using SeDuMi 1.02, a MATLAB toolbox for optimization over symmetric cones. *Optim. Method. Softw.*, 11–12:625–653, 1999.
- [17] R. H. Tütüncü, K. C. Toh and M. J. Todd. Solving semidefinite-quadratic-linear programs using SDPT3 *Math. Program.*, 95:189–217, 2003.
- [18] E. van den Berg and M. P. Friedlander. Probing the Pareto frontier for basis-pursuit solutions. *SIAM J. Sci. Comput.*, 31:890–912, 2008.